
Linear Models of Classification: Probabilistic Generative and Discriminative Models

*Prof. Nicholas Zabaras
Center for Informatics and Computational Science*

<https://cics.nd.edu/>

*University of Notre Dame
Notre Dame, Indiana, USA*

*Email: nzabaras@gmail.com
URL: <https://www.zabaras.com/>*

February 14, 2019

Contents

- ❑ Probabilistic Generative Models for two Classes, Logistic Sigmoid, Models with $K > 2$, Gaussian Class Conditionals
- ❑ Maximum Likelihood Solution, Multiclass Case, Discrete Features
- ❑ Exponential Family, Generalization to Exponential Family Class Conditionals
- ❑ Probabilistic Discriminative Models, Nonlinear Basis in Classification Models, Logistic Regression and Generalized Linear Models, Cross Entropy Error, Sequential Update, Steepest Descent Method, Newton's Method, BFGS, Regularization, Linearly Separable Data, Iterative Reweighted Least Squares
- ❑ Multiclass Logistic Regression

- Following closely Chris Bishop's PRML book, Chapter 4.
- K. Murphy, Machine Learning: A probabilistic Perspective, Chapter 8

Probabilistic Generative Models

- Models with linear decision boundaries arise from simple assumptions about the probabilistic distribution of the data.
- We here adopt **a generative approach**: we model the class-conditional densities $p(\mathbf{x}|C_k)$, as well as the class priors $p(C_k)$, and then use these to compute posterior probabilities $p(C_k|\mathbf{x})$ through Bayes' theorem.
- Let us consider two classes:

$$p(C_1 | \mathbf{x}) = \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_1)p(C_1) + p(\mathbf{x} | C_2)p(C_2)} = \frac{1}{1 + e^{-a}} = \sigma(a)$$

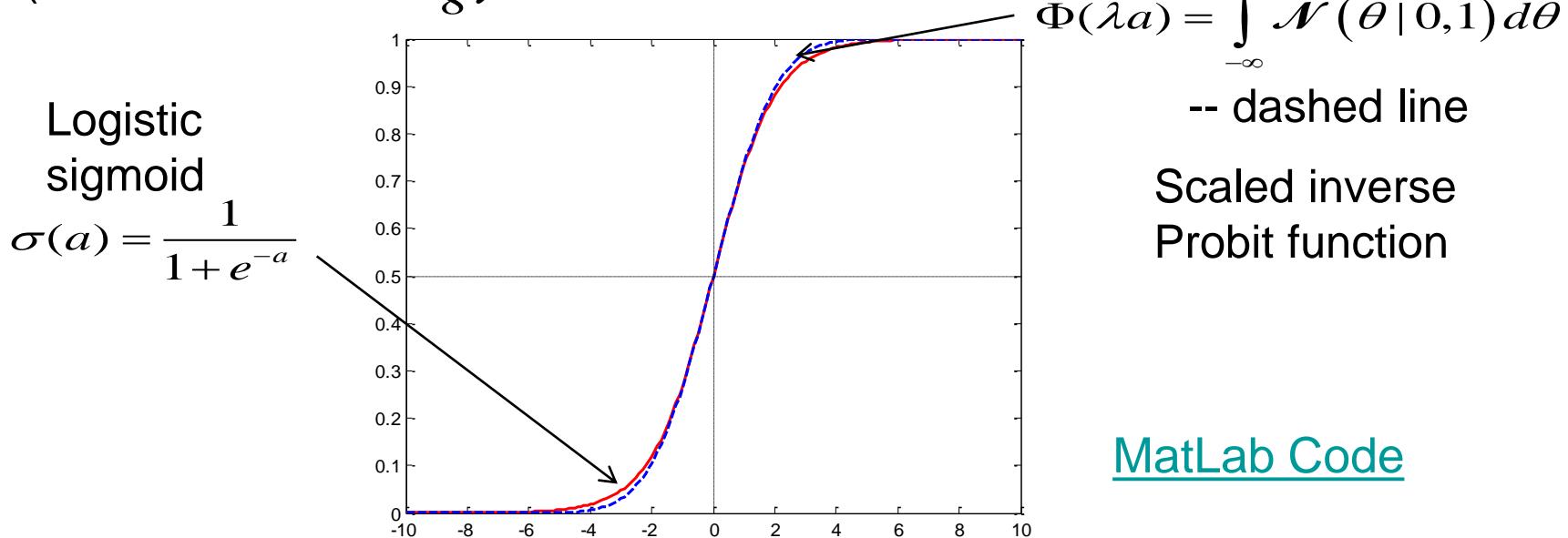
where:

$$a = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} \quad \text{and} \quad \sigma(a) = \frac{1}{1 + e^{-a}}$$

Logistic sigmoid

Logistic Sigmoid (S-shaped)

- We can approximate the logistic sigmoid function with the scaled inverse probit function.
- The derivatives of the two curves are the same for $a = 0$ (selection of $\lambda^2 = \frac{\pi^2}{8}$).



- The logistic sigmoid function satisfies:

$$\sigma(-a) = 1 - \sigma(a)$$

- The *inverse of $\sigma(a)$ is the logit function*:

$$a = \ln \frac{\sigma}{1 - \sigma} = \ln \frac{p(C_1 | \mathbf{x})}{p(C_2 | \mathbf{x})}$$

Probabilistic Generative Models

- The appearance of the logistic sigmoid is rather arbitrary up to this point.

$$p(C_1 | \mathbf{x}) = \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_1)p(C_1) + p(\mathbf{x} | C_2)p(C_2)} = \frac{1}{1+e^{-a}} = \sigma(a)$$

- This will be useful provided $a(\mathbf{x}) = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)}$ takes a simple functional form.
- We will consider problems with $a(\mathbf{x})$ being a linear function of \mathbf{x} .
- In this case, the posterior probability is governed by a generalized linear model.

Probabilistic Generative Models $K > 2$

- For the case of $K > 2$ classes, we have

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{\sum_j p(\mathbf{x} | C_j)p(C_j)} = \frac{e^{a_k}}{\sum_j e^{a_j}}, \quad a_k = \ln(p(\mathbf{x} | C_k)p(C_k))$$

- This is known as the normalized exponential and can be regarded as a multiclass generalization of the logistic sigmoid.
- Here we define:

$$a_k = \ln(p(\mathbf{x} | C_k)p(C_k))$$

- The normalized exponential is known as the softmax function.
- It can be seen as a smoothed version of the max function because, if $a_k \geq a_j$ for all $j \neq k$, then $p(C_k | \mathbf{x}) \cong 1$, and $p(C_j | \mathbf{x}) \cong 0$.

Gaussian Class Conditionals

- Assume that the class-conditional densities are Gaussians and all classes share the same covariance matrix:

$$p(\mathbf{x} | C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}$$

- We consider the case of two classes:

$$p(C_1 | \mathbf{x}) = \sigma(a), a(\mathbf{x}) = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)}$$

where:

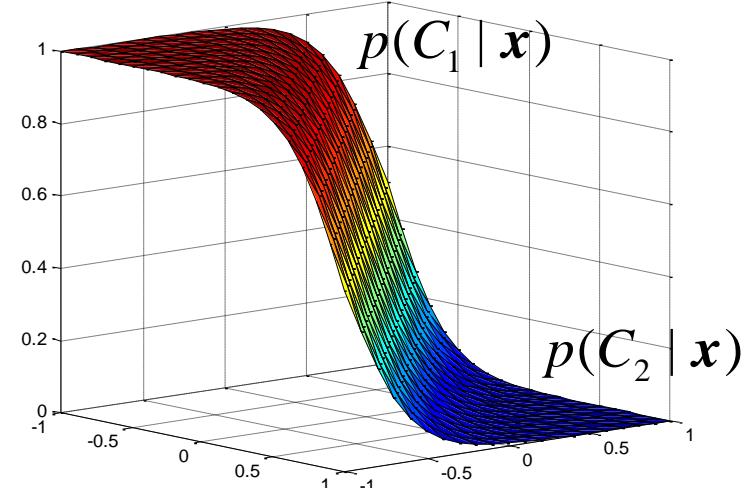
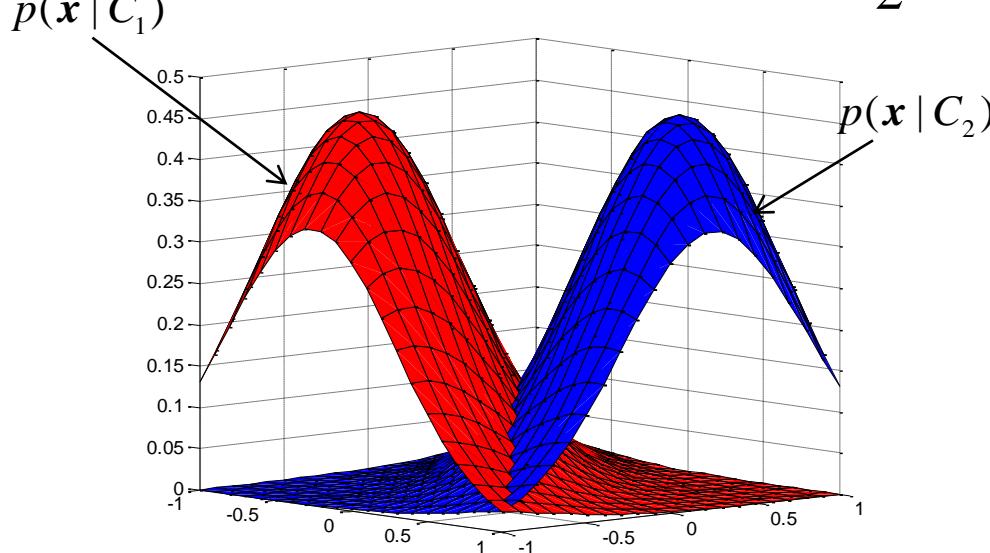
$$\begin{aligned} a(\mathbf{x}) &= -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) + \ln \frac{p(C_1)}{p(C_2)} = \\ &= \underbrace{\left(\boldsymbol{\mu}_1^T \Sigma^{-1} - \boldsymbol{\mu}_2^T \Sigma^{-1} \right) \mathbf{x}}_{w^T} - \underbrace{\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}}_{w_0} \end{aligned}$$

Gaussian Class Conditionals and Posteriors

- We thus obtain a linear function of x :

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

$$\mathbf{w} = \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}$$



[MatLab Code](#)

- Class conditional densities (left) and corresponding posteriors (right).
- On the right, the surface is colored using a proportion of red ink given by $p(C_1 | \mathbf{x})$ and a proportion of blue ink given by $p(C_2 | \mathbf{x}) = 1 - p(C_1 | \mathbf{x})$

Gaussian Class Posteriors

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

- *The decision boundaries correspond to surfaces along which the posterior probabilities $p(C_k | \mathbf{x})$ are constant.*
- They are linear in input space.
- The prior probabilities $p(C_k)$ enter only through the bias parameter w_0

$$w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}$$

so that *changing $p(C_k)$ only leads to parallel shifts of the decision boundaries.*

Probabilistic Generative Models $K > 2$

- For the case of $K > 2$ classes, we have

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{\sum_j p(\mathbf{x} | C_j)p(C_j)} = \frac{e^{a_k}}{\sum_j e^{a_j}} \quad \text{where } a_k = \ln(p(\mathbf{x} | C_k)p(C_k))$$

- The resulting decision boundaries, corresponding to the minimum misclassification rate, occur when the two largest posterior probabilities are equal, and so are defined by linear functions of \mathbf{x} , leading to a generalized linear model.
- For this case, we can see immediately that:

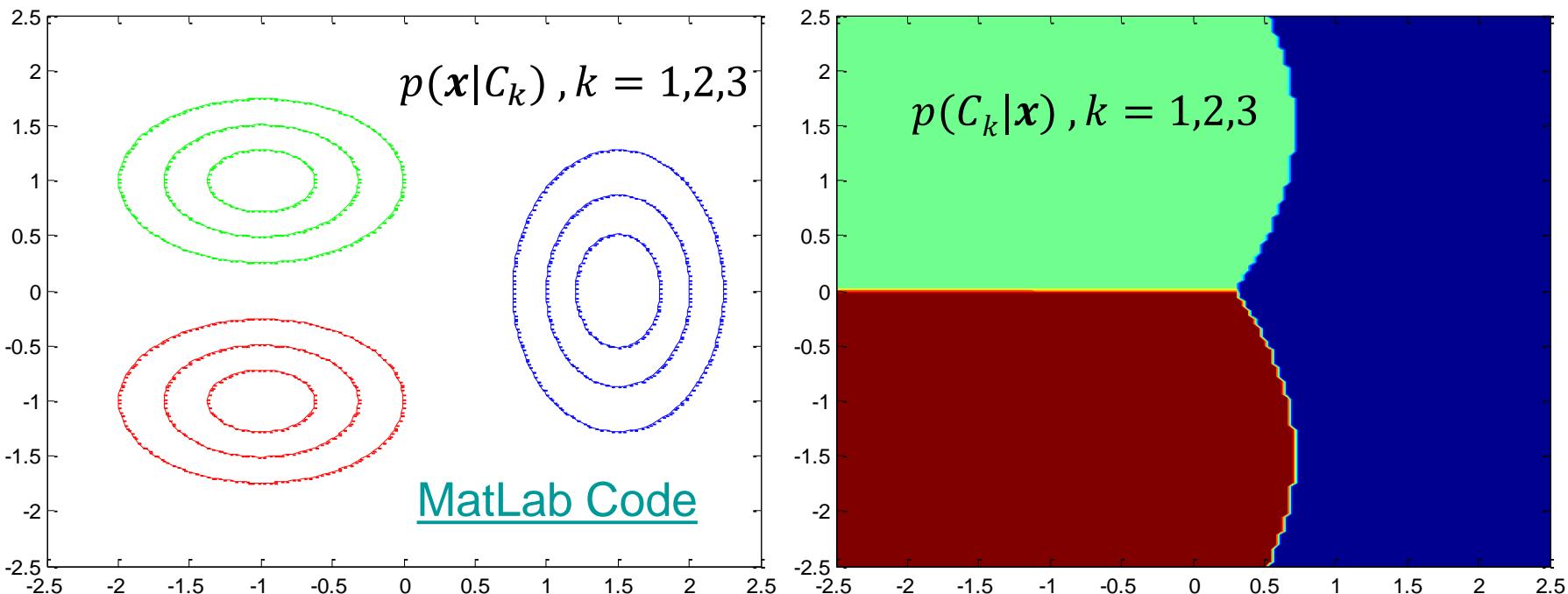
$$a_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$

$$\mathbf{w}_k = \Sigma^{-1} \boldsymbol{\mu}_k, w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln p(C_k)$$

- $a_k(\mathbf{x})$ are linear functions of \mathbf{x} due to the cancellation of the quadratic terms with the shared Σ .

Probabilistic Generative Models $K > 2$

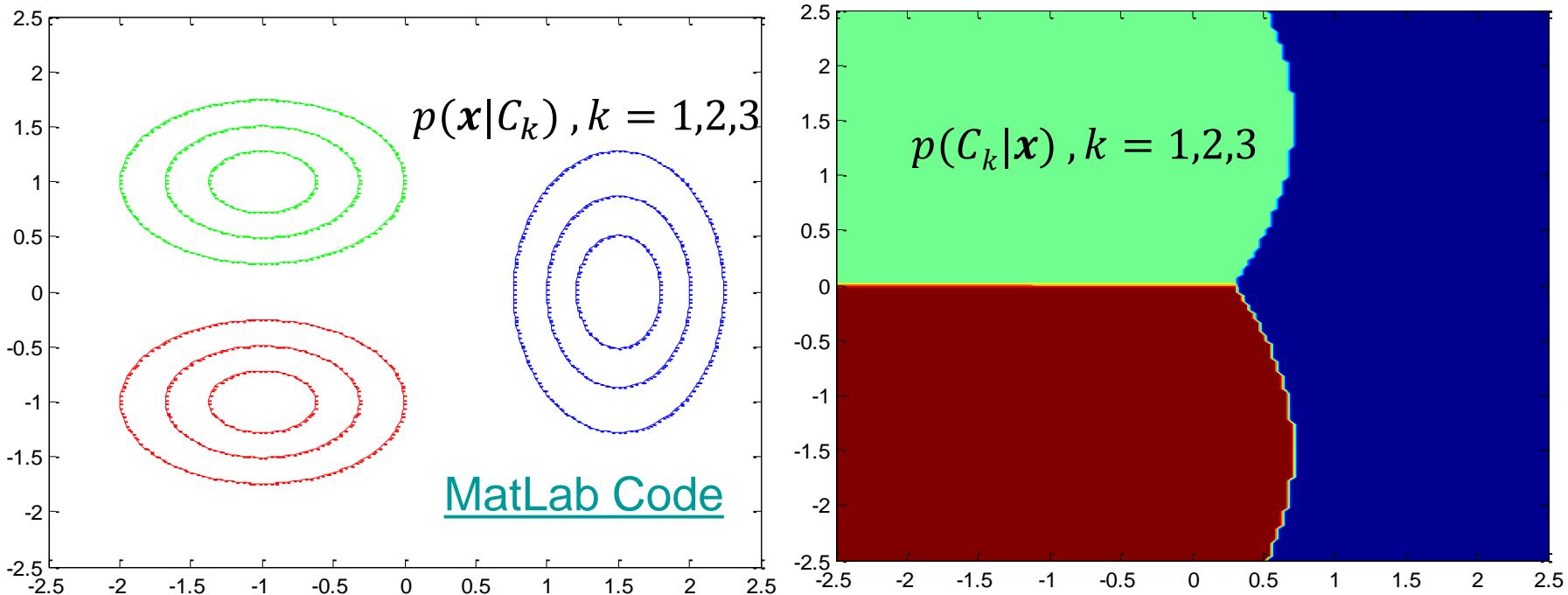
- If we relax the assumption of a shared covariance matrix and allow each $p(x|C_k)$ to have its own Σ_k , then we obtain quadratic functions of x and a quadratic discriminant.



- The left-hand plot shows $p(x|C_k), k = 1,2,3$ each a Gaussian distribution. The red & green classes have the same Σ .

Probabilistic Generative Models $K > 2$

- The right-hand plot shows the posteriors $p(C_k|x)$, $k = 1,2,3$, in which the three colors represent the posterior probabilities for the respective 3 classes (scaled).
- Note that for *the two classes with the same covariance, the decision boundary is linear while for the other 2 is quadratic.*



Maximum Likelihood Solution

- Once we have a parametric form for $p(x|C_k)$, we can compute the parameters, together with the prior class probabilities $p(C_k)$, using MLE.
- This requires a training data set comprising of x and their corresponding class labels.
- Consider the case of two classes each having a Gaussian class-conditional density with a shared covariance matrix.
- Suppose we have a data set $\{x_n, t_n\}$ where $n = 1, \dots, N$. Here $t_n = 1$ denotes class C_1 and $t_n = 0$ denotes C_2 .
- We denote the prior class probabilities $p(C_1) = \pi$, $p(C_2) = 1 - \pi$.

Maximum Likelihood Solution

- For a data point x_n from class C_1 , we have $t_n = 1$ and hence

$$p(x_n, C_1) = p(C_1)p(x_n | C_1) = \pi \mathcal{N}(x_n | \mu_1, \Sigma)$$

- Similarly for class C_2 , we have $t_n = 0$

$$p(x_n, C_2) = p(C_2)p(x_n | C_2) = (1 - \pi) \mathcal{N}(x_n | \mu_2, \Sigma)$$

- Thus *the likelihood function is given by*

$$p(t, X / \pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N [\pi \mathcal{N}(x_n | \mu_1, \Sigma)]^{t_n} [(1 - \pi) \mathcal{N}(x_n | \mu_2, \Sigma)]^{1-t_n}$$

where $t = (t_1, \dots, t_N)^T$.

- We maximize the log likelihood function.

Maximum Likelihood Solution

$$p(\mathbf{t}, \mathbf{X} / \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1-\pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n}$$

- Consider first the maximization with respect to π . The terms of interest in the log-likelihood are:

$$\sum_{n=1}^N \{t_n \ln \pi + (1-t_n) \ln(1-\pi)\}$$

- Taking derivative wrt π and setting it equal to zero:

$$\pi = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N_1 + N_2}$$

where N_i the total number of data points in class C_i .

- As expected, the maximum likelihood estimate for π is simply the fraction of points in class C_1 .
- This result is easily generalized to the multiclass case (see next) where the MLE of the prior probability associated with class C_k is given by the fraction of the training set points N_k .

Maximum Likelihood Solution: Multiclass

$$p(\{\phi_n, t_n\} / \{\pi_k\}) = \prod_{n=1}^N \prod_{k=1}^K [\pi_k p(\phi_n | C_k)]^{t_{nk}}$$

- For the multiclass case, the log likelihood takes the form:

$$\ln p(\{\phi_n, t_n\} / \{\pi_k\}) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} (\ln \pi_k + \ln p(\phi_n | C_k)), \quad \sum_{k=1}^K \pi_k = 1$$

- To maximize the log likelihood, we consider:

$$\sum_{n=1}^N \sum_{k=1}^K t_{nk} (\ln \pi_k + \ln p(\phi_n | C_k)) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

- Taking the derivative wrt π_k we arrive at the same result as for two classes:

$$\sum_{n=1}^N \frac{t_{nk}}{\pi_k} + \lambda = 0 \Rightarrow \pi_k \lambda = -\sum_{n=1}^N t_{nk} = -N_k \Rightarrow \begin{cases} \sum_{k=1}^K \pi_k \lambda = -N \\ \pi_k \lambda = -N_k \end{cases} \Rightarrow \begin{cases} \lambda = -N \\ \pi_k \lambda = -N_k \end{cases} \Rightarrow \pi_k = \frac{N_k}{N}$$

Maximum Likelihood Solution

$$p(\mathbf{t}, \mathbf{X} | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_{n=1}^N [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1-\pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n}$$

- Consider now the maximization with respect to $\boldsymbol{\mu}_1$. The terms of interest in the log-likelihood are:

$$\sum_{n=1}^N t_n \ln \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) = -\frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) + const.$$

- Taking derivative wrt $\boldsymbol{\mu}_1$ and setting it equal to zero:

$$\boldsymbol{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n, \text{ and similarly } \boldsymbol{\mu}_2 = \frac{1}{N_2} \sum_{n=1}^N (1-t_n) \mathbf{x}_n$$

- Thus the MLE for $\boldsymbol{\mu}_i$ is simply the mean of the fraction of points in class C_i .
- This result is easily generalized to the multiclass case. The MLE of the prior probability associated with class C_k is given by the fraction of the training set points N_k .

Maximum Likelihood Solution

$$p(t | \pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N [\pi \mathcal{N}(x_n | \mu_1, \Sigma)]^{t_n} [(1-\pi) \mathcal{N}(x_n | \mu_2, \Sigma)]^{1-t_n}$$

- Consider finally the maximization with respect to Σ . The terms of interest in the log-likelihood are:

$$\begin{aligned} & -\frac{1}{2} \sum_{n=1}^N t_n \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N t_n (\mathbf{x}_n - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_1) \\ & - \frac{1}{2} \sum_{n=1}^N (1-t_n) \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (1-t_n) (\mathbf{x}_n - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_2) \\ & = -\frac{N}{2} \ln |\Sigma| - \frac{N}{2} \text{Tr}(\Sigma^{-1} S), \text{ where } S = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2 \end{aligned}$$

$$S_1 = \frac{1}{N_1} \sum_{n \in C_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)(\mathbf{x}_n - \boldsymbol{\mu}_1)^T, S_2 = \frac{1}{N_2} \sum_{n \in C_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)(\mathbf{x}_n - \boldsymbol{\mu}_2)^T$$

- Using the result for the MLE for a Gaussian distribution, we obtain $\Sigma = S$. This result is a weighted average of the covariance matrices associated with each class separately.

Maximum Likelihood Solution

$$\pi = \frac{N_1}{N_1 + N_2}, \quad \boldsymbol{\mu}_1 = \frac{1}{N_1} \sum_{n=1}^N t_n \mathbf{x}_n, \quad \boldsymbol{\mu}_2 = \frac{1}{N_2} \sum_{n=1}^N (1-t_n) \mathbf{x}_n$$

$$\Sigma = \frac{N_1}{N} \mathbf{S}_1 + \frac{N_2}{N} \mathbf{S}_2, \quad \mathbf{S}_1 = \frac{1}{N_1} \sum_{n \in C_1} (\mathbf{x}_n - \boldsymbol{\mu}_1)(\mathbf{x}_n - \boldsymbol{\mu}_1)^T, \quad \mathbf{S}_2 = \frac{1}{N_2} \sum_{n \in C_2} (\mathbf{x}_n - \boldsymbol{\mu}_2)(\mathbf{x}_n - \boldsymbol{\mu}_2)^T$$

- These results are easily extended to the K class problem to obtain the corresponding MLE solutions for the parameters (for Gaussian class-conditionals with a shared covariance matrix).
- As we have seen before, **MLE is not robust to outliers.**
 - Fitting Gaussian distributions to the classes is not robust to outliers, because the MLE of a Gaussian is not robust.

Discrete Features: Naive Bayes Assumption

- Here we consider the case of discrete feature values x_i . Consider binary feature values $x_i \in \{0, 1\}, i = 1, \dots, D$.
- For D inputs, a general distribution corresponds to a table of 2^D numbers for each class, containing $2^D - 1$ independent variables (due to the summation constraint).
- This grows exponentially with the number D of features.
- Consider the **Naive Bayes Assumption**: the feature values are treated as independent, conditioned on the class C_k .
- Thus we have class-conditional distributions of the form

$$p(\mathbf{x} | C_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i}$$

Discrete Features

$$p(\mathbf{x} | C_k) = \prod_{i=1}^D \mu_{ki}^{x_i} (1 - \mu_{ki})^{1-x_i}$$

- There are D independent parameters for each class.
- Here μ_{ki} is the probability the i^{th} feature takes the value 1.
- Substituting into $a_k = \ln(p(\mathbf{x} | C_k)p(C_k))$ gives:

$$a_k(\mathbf{x}) = \sum_{i=1}^D \left\{ x_i \ln \mu_{ki} + (1 - x_i) \ln (1 - \mu_{ki}) \right\} + \ln p(C_k)$$

- These are linear functions of the input values x_i .
 - For $K = 2$, we can alternatively consider the logistic sigmoid
- $$p(C_1 | \mathbf{x}) = \sigma(a), a(\mathbf{x}) = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)}$$
- Analogous results are obtained for discrete variables each of which can take $M > 2$ states.

Exponential Family

- We have seen that for both Gaussian and discrete inputs, the posterior class probabilities are given by generalized linear models with logistic sigmoid ($K = 2$) or softmax ($K \geq 2$ classes) activation functions.
- This result is general for the class-conditional densities $p(x|C_k)$ that are members of the exponential family.
- For the exponential family, the distribution of x can be written in the form

$$p(x | \lambda_k) = h(x)g(\lambda_k)\exp\left\{\lambda_k^T u(x)\right\}$$

- Here, we consider the case $u(x) = x$.

Exponential Family

$$p(\mathbf{x} | \boldsymbol{\lambda}_k) = h(\mathbf{x})g(\boldsymbol{\lambda}_k)\exp\left\{\boldsymbol{\lambda}_k^T \mathbf{u}(\mathbf{x})\right\}$$

- With $\mathbf{u}(\mathbf{x}) = \mathbf{x}$ and by introducing a scaling parameter s as $p(\mathbf{x} | s) = (1/s)f(\mathbf{x}/s)$, we obtain the restricted set of exponential family class-conditional densities of the form:

$$p(\mathbf{x} | \boldsymbol{\lambda}_k, s) = \frac{1}{s} h\left(\frac{1}{s}\mathbf{x}\right) g(\boldsymbol{\lambda}_k) \exp\left\{\frac{1}{s} \boldsymbol{\lambda}_k^T \mathbf{x}\right\}$$

- Here each class has its own parameter $\boldsymbol{\lambda}_k$ but all classes share the same scale parameter s .
- For the 2 class problem, we substitute this expression for the class-conditional densities into

$$a(\mathbf{x}) = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)}$$

Exponential Family

$$p(\mathbf{x} | \boldsymbol{\lambda}_k, s) = \frac{1}{s} h\left(\frac{1}{s} \mathbf{x}\right) g(\boldsymbol{\lambda}_k) \exp\left\{\frac{1}{s} \boldsymbol{\lambda}_k^T \mathbf{x}\right\} \quad a(\mathbf{x}) = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)}$$

- The posterior class probability is again given by a logistic sigmoid acting on a linear function $a(\mathbf{x})$

$$p(C_1 | \mathbf{x}) = \sigma(a(\mathbf{x}))$$

which is given by

$$a(\mathbf{x}) = \frac{1}{s} (\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2)^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_1) - \ln g(\boldsymbol{\lambda}_2) + \ln p(C_1) - \ln p(C_2)$$

- Similarly for the K –class problem, we substitute the class-conditionals directly in $a_k(\mathbf{x}) = \ln(p(\mathbf{x} | C_k)p(C_k))$ to obtain:

$$a_k(\mathbf{x}) = \frac{1}{s} \boldsymbol{\lambda}_k^T \mathbf{x} + \ln g(\boldsymbol{\lambda}_k) + \ln p(C_k)$$

which is linear in \mathbf{x} .

Probabilistic Discriminative Models

- For the two-class problem, the posterior probability of class C_k can be written as a logistic sigmoid acting on a linear function of x , for a wide choice of $p(x|C_k)$.
- For the multiclass problem, the posterior probability of C_k is given by a softmax transformation of a linear function of x .
- For specific choices of $p(x|C_k)$, we used MLE to determine the parameters of the densities and the class priors $p(C_k)$ and then used Bayes' theorem to find the posterior class probabilities.
- An alternative approach that we follow in this lecture is to *use the functional form of the generalized linear model explicitly and determine its parameters directly by using MLE*.
- There is an efficient algorithm finding such solutions known as iterative reweighted least squares, or IRLS.

Probabilistic Discriminative Models

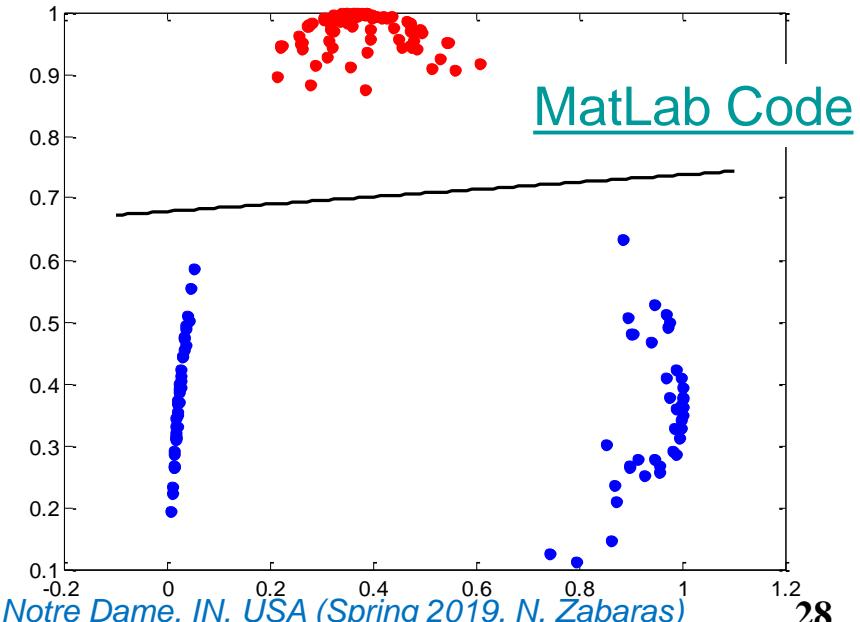
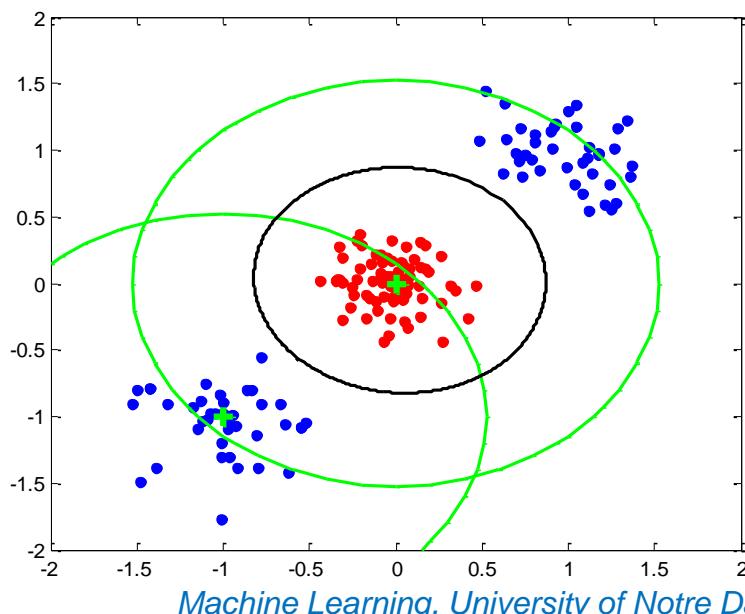
- Indirect approach – Generative Modeling: find the parameters of the generalized linear model, by fitting class-conditional densities and class priors separately and then applying Bayes' theorem.
 - With this model, we *can generate synthetic data by drawing values of x from the marginal $p(x)$.*
- Direct approach – Discriminative Modeling: we are maximizing a likelihood function defined through $p(C_k|x)$, which represents a form of discriminative training.
- An advantage of the discriminative approach is that there are typically fewer adaptive parameters to be determined.
- It leads to improved predictive performance when the class-conditional density assumptions are a poor approximation to the true distributions.

Fixed Basis Functions

- First make a fixed nonlinear transformation of the inputs using a vector of basis functions $\phi(x)$, then apply the algorithms discussed up to now.
- The resulting decision boundaries will be linear in the feature space ϕ , and these correspond to nonlinear decision boundaries in the original x space.
- Classes that are linearly separable in the feature space $\phi(x)$ need not be linearly separable in the original input space x .

Nonlinear Basis in Classification Models

- Left: original space (x_1, x_2) with data from red & blue classes
- 2 ‘Gaussian’ basis functions $\phi_1(x), \phi_2(x)$ are defined in (x_1, x_2) with green crosses the centers & green circles the contours.
- Right: the corresponding feature space (ϕ_1, ϕ_2) together with the linear decision boundary obtained by a logistic regression model. This corresponds to a nonlinear decision boundary in (x_1, x_2) (black curve on the left).



Fixed Basis Functions

- One of the basis functions is typically set to a constant, say $\phi_0(x) = 1$, so that w_0 plays the role of a bias.
- For problems with significant overlap between the class-conditional densities $p(x|C_k)$, the posterior probabilities $p(C_k|x)$ at least for some values of x may not be 0 or 1.
 - In such cases, the optimal solution is obtained by modeling the posterior probabilities accurately and then applying standard decision theory.
- The transformation $\phi(x)$ cannot remove class overlap, it can increase the level of overlap, or create overlap where none existed in the x space.
- However, **suitable choices of nonlinearity** can make the process of modeling the posterior probabilities easier.

Fixed Basis Functions

- Fixed basis function models have serious limitations.
- We should allow the basis functions to adapt to the data.
- Still models with fixed nonlinear basis functions are important as they introduce many of the key concepts needed for an understanding models with data-adaptive basis.

Logistic Regression Generalized Linear Models

- Consider a simple generalized linear problem with a two-class classification.
- In generative models we saw that the posterior probability of class C_1 can be written as a logistic sigmoid acting on a linear function of the feature vector ϕ so that

$$p(C_1|\phi) = y(\phi) = \sigma(w^T\phi)$$

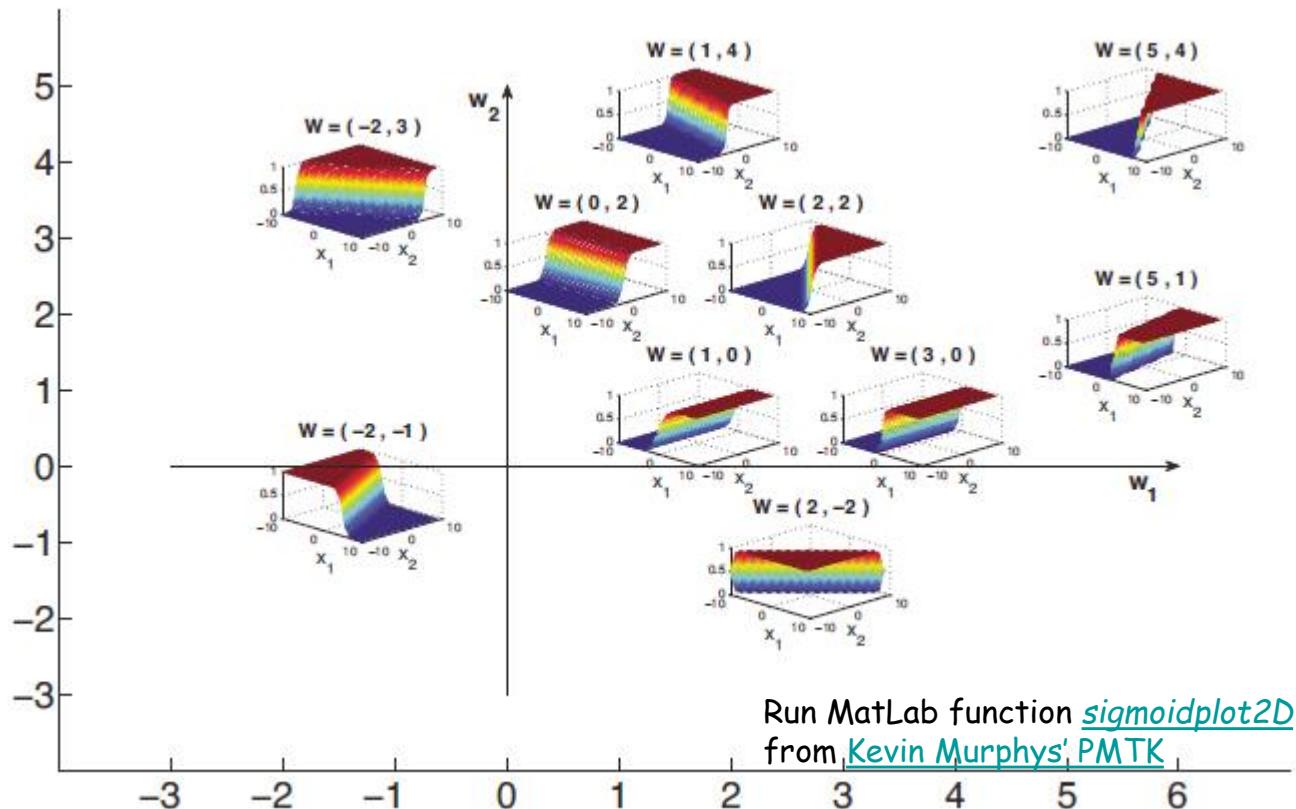
with $p(C_2|\phi) = 1 - p(C_1|\phi)$. Here $\sigma(\cdot)$ is the logistic sigmoid function defined by

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

- This model is known as logistic regression – but *note its a model for classification not regression.*

Plots of $p(C_1|f) = y(\phi) = \sigma(w^T\phi)$

- Plots of $p(C_1|f) = \sigma(w_1x_1 + w_2x_2)$ for a 2D input and different w .
- If we threshold these probabilities at 0.5, we induce a linear decision boundary whose normal is w .
- To the right of this have $\sigma(w^T x) > 0$ and to the left $\sigma(w^T x) < 0$.



Logistic Regression

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$

- For an M –dimensional feature space ϕ , this model has M adjustable parameters.
- If we had fitted Gaussian class conditional densities using MLE, we would have used $2M$ parameters for the means and $M(M + 1)/2$ parameters for the (shared) covariance. Together with the class prior $p(C_1)$, this gives a total of $M(M + 5)/2 + 1$ parameters, i.e. $\mathcal{O}(M^2)$, in contrast to $\mathcal{O}(M)$ dependence on M in logistic regression.
- For large M , there is advantage in working with logistic regression.
- We now use MLE to determine the parameters in logistic regression. We will use:

$$\frac{d\sigma}{da} = \sigma(1 - \sigma)$$

Logistic Regression - Cross-Entropy Error

- For a data set $\{\phi_n, t_n\}$, where $t_n \in \{0, 1\}$ and $\phi_n = \phi(x_n)$, with $n = 1, \dots, N$, the likelihood function can be written

$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}, \mathbf{t} = (t_1, \dots, t_N)^T, y_n = p(C_1 | \phi_n)$$

- Using this, we define the cross-entropy error function as:

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1 - t_n) \ln (1 - y_n) \right\}, y_n = \sigma(a_n), a_n = \mathbf{w}^T \phi_n$$

- Taking the gradient wrt \mathbf{w} and using $\frac{d\sigma}{da} = \sigma(1 - \sigma)$ gives:

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n, y_n = \sigma(a_n), a_n = \mathbf{w}^T \phi_n$$

- The contribution to $\nabla E(\mathbf{w})$ from data point n is given by the ‘error’ $y_n - t_n$ between the target value and the prediction of the model, times the basis ϕ_n . This is identical to the gradient of the sum-of-squares error function for linear regression.

Logistic Regression - Sequential Update

- We use this result

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_n, \quad y_n = \sigma(a_n), \quad a_n = \mathbf{w}^T \boldsymbol{\phi}_n$$

to derive a sequential algorithm in which data are presented one at a time.

- The weight vector is updated using

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n = \mathbf{w}^{(\tau)} - \eta (y_n - t_n) \boldsymbol{\phi}_n$$

- Note that for linearly separable data (the hyperplane $\sigma = 0.5$, $\mathbf{w}^T \boldsymbol{\phi} = 0$ separates the two classes) any decision boundary separating the two classes needs to satisfy:

$$\mathbf{w}^T \boldsymbol{\phi}_n \begin{cases} \geq 0 & \text{if } t_n = 1 \\ < 0 & \text{otherwise} \end{cases}$$

- However we can see from the top Eq. that the negative log likelihood is minimized when: $y_n = \sigma(\mathbf{w}^T \boldsymbol{\phi}_n) = t_n$, $n = 1, \dots, N$. This requires that the sigmoid function saturates which occurs when its argument

$$\mathbf{w}^T \boldsymbol{\phi}_n \rightarrow \pm\infty, \text{ i.e. when } \|\mathbf{w}\| \rightarrow \infty.$$

Logistic Regression Linearly Separable Data

- MLE can exhibit severe over-fitting for data sets *that are linearly separable*.
- Indeed for linearly separable data, the hyperplane corresponding to $P(C_1|x) = \sigma = 0.5 \rightarrow \mathbf{w}^T \phi = 0$ separates the two classes and will have the following:

$$\mathbf{w}^T \phi_n = \begin{cases} \geq 0 & \text{if } t_n = 1 \\ < 0 & \text{if } t_n = 0 \end{cases}$$

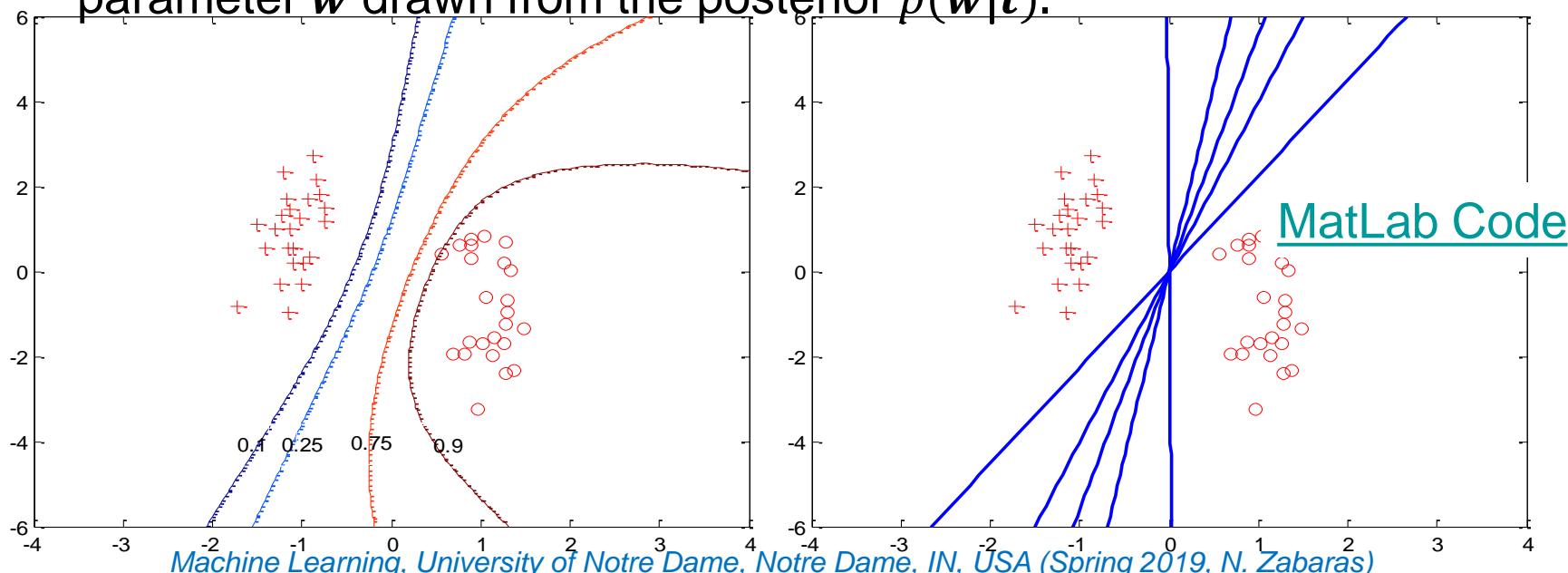
- In addition we can see that the negative log-likelihood is minimized when
- $$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = 0 \Rightarrow y_n = \sigma(\mathbf{w}^T \phi_n) = t_n \quad \forall n, \text{i.e. } y_n = t_n = 0 \text{ or } y_n = t_n = 1$$
- This occurs when *the sigmoid function is saturated*, which occurs when its argument $\mathbf{w}^T \phi$ goes to $\pm\infty$, i.e. *when the magnitude of w goes to infinity*.
 - The logistic sigmoid function becomes infinitely steep in feature space (step function) so every training point from each class k is assigned a posterior probability $p(C_k | x) = 1$.
 - There are infinite solutions since any separating hyperplane will give the same posterior probabilities at the training data points.

Logistic Regression Linearly Separable Data

- Maximum likelihood provides no way to favor one such solution over another.
- *Which solution is found depends on the choice of optimization algorithm and on the parameter initialization.*
- The problem arises even if the number of data points is large compared to the number of parameters in the model, so long as the training data set is linearly separable.
- The singularity can be avoided by inclusion of a prior and finding a MAP solution for w , or equivalently by adding a regularization term to the error function.

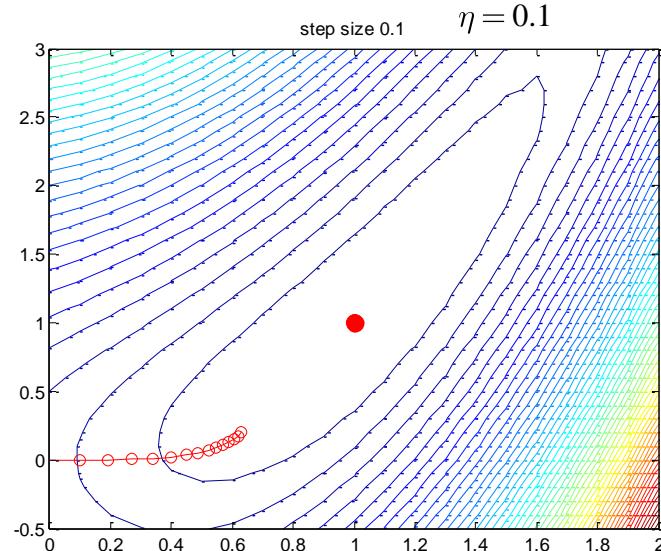
Logistic Regression Linearly Separable Data

- Bayesian approach to logistic regression for linearly separable data set.
- Left: the predictive distribution using Variational Inference (details will be given in another lecture). The decision boundary lies between the clusters of data points. The contours of the predictive distribution splay out away from the data reflecting the greater uncertainty in the classification of such regions.
- Right: the decision boundaries corresponding to 5 samples of the parameter w drawn from the posterior $p(w|t)$.



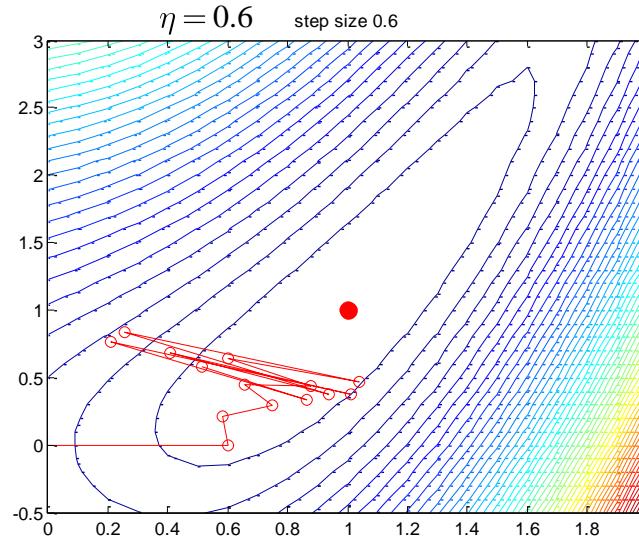
Gradient Descent Example

- The effect of the learning rate η on the performance of a steepest descent algorithm is shown on a simple function.



$$f(\theta) = 0.5 \theta_1^2 - \theta_2^2 + 0.5 \theta_1 - 1^2$$

$$\theta_{k+1} = \theta_k - \eta_k \mathbf{g}_k = \theta_k - \eta_k \nabla f(\theta)|_k$$

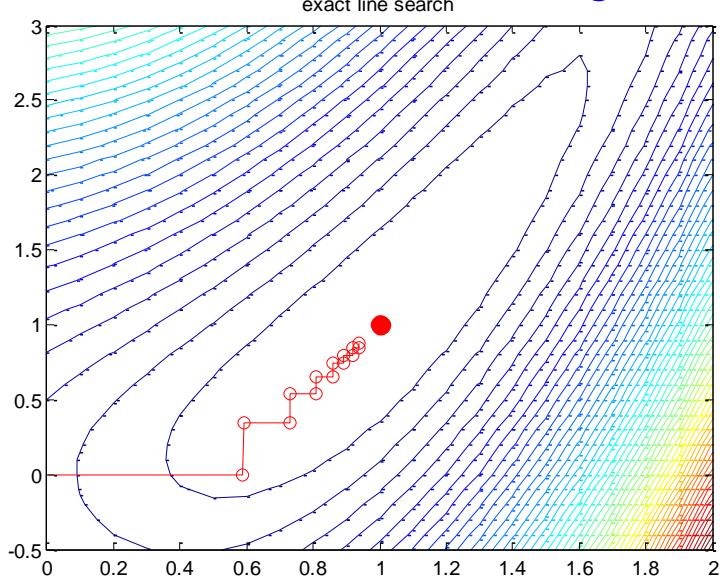


[steepestDescentDemo](#)
from [Kevin Murphys' PMTK](#)

- Small step size leads to slow convergence and large step size to never converging. In both cases, the step size is kept constant.

Gradient Descent Example

- Gradient descent algorithm with line search is shown below.



$$f(\theta) = 0.5 \theta_1^2 - \theta_2^2 + 0.5 \theta_1 - 1^2$$

$$\theta_{k+1} = \theta_k - \eta_k g_k, g = \nabla f$$

$$\text{Pick } \eta \text{ to minimize : } \phi(\eta) = f(\theta_k + \eta d_k) \approx f(\theta_k) + \eta \nabla f(\theta_k)^\top d_k$$

[steepestDescentDemo](#)

from [Kevin Murphys' PMTK](#)

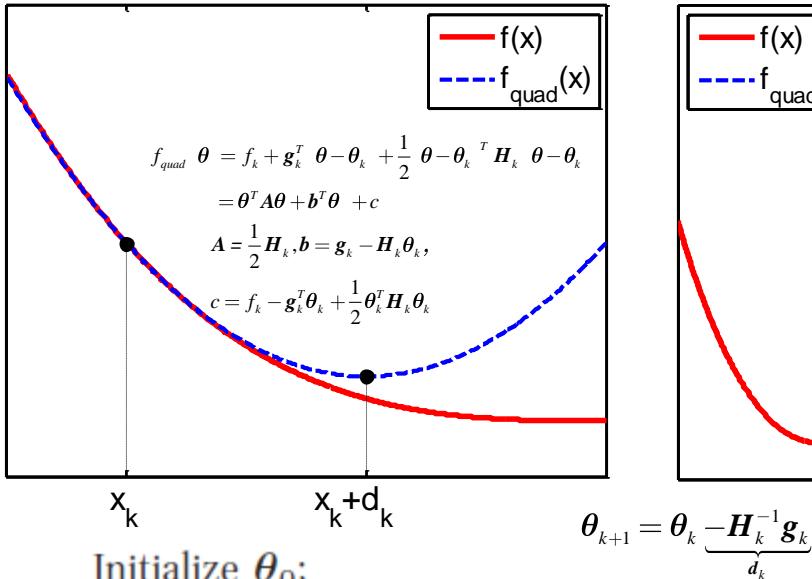
- Bertsekas, D. (1999). [Nonlinear Programming](#) (Second ed.). Athena Scientific.
- Nocedal, J. and S. Wright (2006). [Numerical Optimization](#). Springer.
- Golub, G. and C. F. van Loan (1996). [Matrix computations](#). Johns Hopkins University Press.



- Consecutive directions are orthogonal: noticing that $\phi'(\eta) = 0$ leads to: $d^T g_{k+1} = 0, g_{k+1} = f'(\theta_k + \eta d_k)$. So at the end of each step k, either $\nabla f|_{k+1} = 0$ (stationary point was achieved) or $d \perp \nabla f|_{k+1}$ (i.e. the search stops at a point where the gradient is normal to the search direction).
- To avoid the zig-zag path, modify as: $\theta_{k+1} = \theta_k - \eta_k g_k + \mu_k (\theta_k - \theta_{k-1})$, $\mu_k \in [0, 1]$. This is the so called *heavy-ball method*.
- Alternatively, one can apply the *conjugate gradient method*.

Newton's Method

- Newton's method takes the Hessian into account: $\theta_{k+1} = \theta_k - \eta_k \mathbf{H}_k^{-1} \mathbf{g}_k$, $\mathbf{g} = \nabla f$



Initialize θ_0 ;

for $k = 1, 2, \dots$ until convergence do

Evaluate $\mathbf{g}_k = \nabla f(\theta_k)$;

Evaluate $\mathbf{H}_k = \nabla^2 f(\theta_k)$;

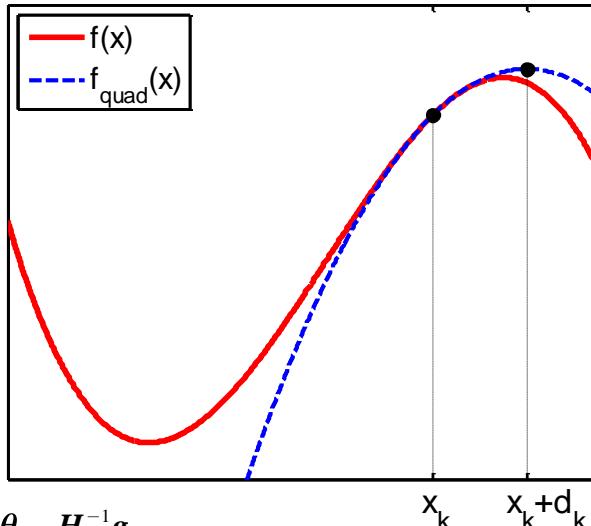
Solve $\mathbf{H}_k \mathbf{d}_k = -\mathbf{g}_k$ for \mathbf{d}_k ;

Use line search to find stepsize η_k along \mathbf{d}_k ;

$\theta_{k+1} = \theta_k + \eta_k \mathbf{d}_k$;

- Vandenberghe, L. (2006). [Applied numerical computing](#): Lecture notes.

- Also one can use the *truncated Newton method* (solve $\mathbf{H}_k \mathbf{d}_k = -\mathbf{g}_k$ with CG but truncate the iterations appropriately)



[newtonsMethodMinQuad](#)
 and [newtonsMethodNonConvex](#)
 from [Kevin Murphys' PMTK](#)

- The step size $d_k = -\mathbf{H}_k^{-1} \mathbf{g}_k$ is what should be added to θ_k to minimize the 2nd order approximation of f around θ_k .

$$f_{quad}(\theta) = 0 \Rightarrow A\theta = -\frac{1}{2}b \Rightarrow \mathbf{H}_k \theta_{k+1} = -\mathbf{g}_k - \mathbf{H}_k \theta_k \\ \Rightarrow \theta_{k+1} = \theta_k - \mathbf{H}_k^{-1} \mathbf{g}_k$$

- H needs to be pos. definite (holds for convex functions).

- If not, d_k may not be a descent direction (Fig. on right).
- For such cases, use the *Levenberg Marquardt method* (adaptively switch between Newton and Steepest descent updates $d_k = -\mathbf{g}_k$)

Limited Memory BFGS

- Computing \mathbf{H} is expensive. Quasi Newton methods iteratively build approximations using the gradient vector at each step.
- The BFGS method builds an approximation \mathbf{B} (rank –two updates) to the Hessian as (we start with $\mathbf{B}_0 = \mathbf{I}$):

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{\mathbf{y}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} - \frac{\mathbf{B}_k \mathbf{s}_k \mathbf{B}_k \mathbf{s}_k^T}{\mathbf{s}_k^T \mathbf{B}_k \mathbf{s}_k}$$

$$\mathbf{s}_k = \boldsymbol{\theta}_k - \boldsymbol{\theta}_{k-1}$$

$$\mathbf{y}_k = \mathbf{g}_k - \mathbf{g}_{k-1}$$

- Often we provide updates to the inverse of the Hessian:

$$\mathbf{C}_{k+1} = \left(\mathbf{I} - \frac{\mathbf{s}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} \right) \mathbf{C}_k \left(\mathbf{I} - \frac{\mathbf{s}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} \right) + \frac{\mathbf{s}_k \mathbf{s}_k^T}{\mathbf{y}_k^T \mathbf{s}_k}, \mathbf{C}_k \approx \mathbf{H}_k^{-1}$$

- To avoid storing the Hessian or its inverse, *limited memory BFGS (L-BFGS)* is used where $\mathbf{H}_k^{-1} \mathbf{g}_k$ is performed by a sequence of inner products with \mathbf{s}_k and \mathbf{y}_k using only the m most recent (~ 20) pairs of $(\mathbf{s}_k, \mathbf{y}_k)$. This requires only $\mathcal{O}(mD)$ storage vs $\mathcal{O}(D^2)$ for the full Hessian.

- Nocedal, J. and S. Wright (2006). *Numerical Optimization*. Springer.

L_2 Regularization

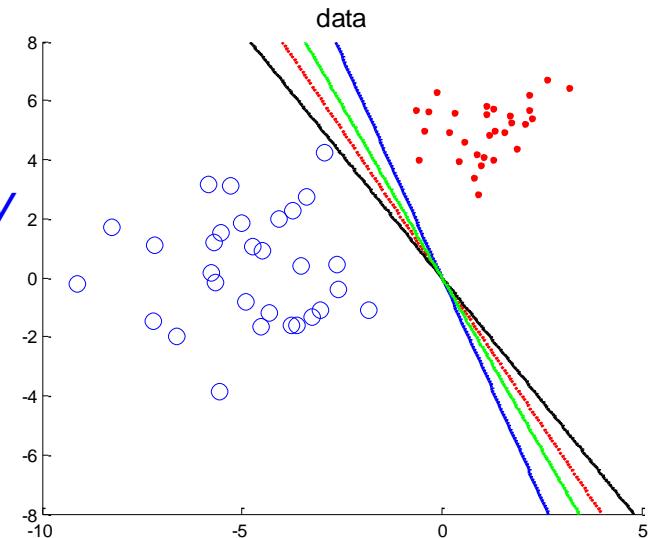
- ❑ Regularization is important for classification even with lots of data.
- ❑ For linearly separable data, the MLE is obtained by $\|\mathbf{w}\| \rightarrow \infty$, corresponding to an infinitely steep sigmoid. In the two class linearly separable data, we show the Log-likelihood.
- ❑ The line is drawn from the origin in the direction of the MLE (which is at ∞). The numbers correspond to the lines on the top Fig.
- ❑ This solution is very brittle and does not generalize well.
- ❑ To prevent this, we use L_2 regularization.

$$f(\mathbf{w}) = NLL(\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w} \text{ (objective function)}$$

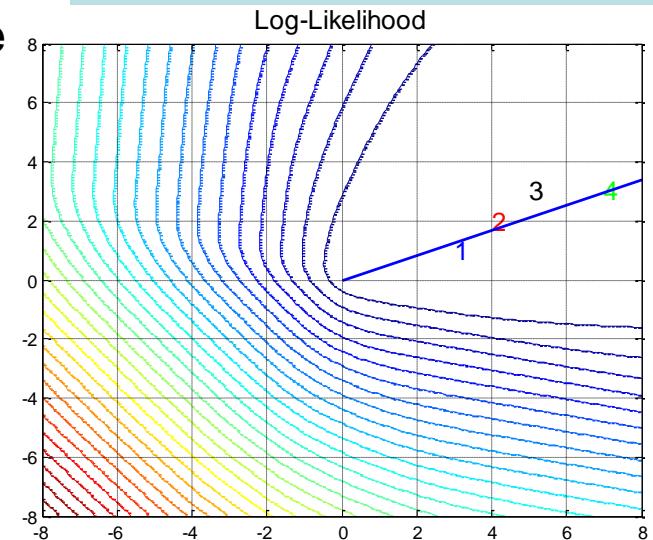
$$\mathbf{g}(\mathbf{w}) = \mathbf{g}(\mathbf{w}) + \lambda \mathbf{w} \text{ (gradient)}$$

$$\mathbf{H}(\mathbf{w}) = \mathbf{H}(\mathbf{w}) + \lambda \mathbf{I} \text{ (Hessian)}$$

- ❑ These modifications provide no difficulty with any gradient-based optimization toolbox.



[logregLaplaceGirolamiDemo](#)
from [Kevin Murphys' PMTK](#)



Iterative Reweighted Least Squares

- In linear regression models, the MLE for a Gaussian noise model had a closed-form solution. This was a consequence of the quadratic dependence of the log likelihood function on w .
- For logistic regression, there is no closed-form solution due to the nonlinearity of the logistic sigmoid function.
- However, the error function is convex and hence has a unique minimum.
- In addition, it can be minimized by an iterative technique based on *Newton-Raphson, using a quadratic approximation to the log likelihood function.*
- The Newton-Raphson update, for minimizing a function $E(w)$, takes the form

$$w^{(new)} = w^{(old)} - H^{-1} \nabla E(w), \quad H_{ij} = \partial^2 E(w) / \partial w_i \partial w_j \quad (\text{Hessian})$$

Iterative Reweighted Least Squares

- Let us first apply the Newton-Raphson method to the linear regression model with the sum-of-squares error function. The gradient and Hessian of this error function are given by

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \phi(x_n))^2 \Rightarrow \nabla E(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \phi(x_n) - t_n) \phi_n = \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{t}$$

where Φ is the $N \times M$ matrix whose n^{th} row is given by ϕ_n^T .

- The Newton-Raphson update takes the form:

$$\mathbf{w}^{(new)} = \mathbf{w}^{(old)} - (\Phi^T \Phi)^{-1} (\Phi^T \Phi \mathbf{w}^{(old)} - \Phi^T \mathbf{t}) = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

where $\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N \phi_n \phi_n^T = \Phi^T \Phi$

- This is the standard least-squares solution to *linear regression*. The error function is quadratic and hence the Newton-Raphson gives the exact solution in one step!

Iterative Reweighted Least Squares

- We similarly apply the Newton-Raphson to the cross-entropy error function for the logistic regression model:

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1-t_n) \ln (1-y_n) \right\}, \quad y_n = \sigma(a_n), \quad a_n = \mathbf{w}^T \boldsymbol{\phi}_n \Rightarrow$$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_n = \boldsymbol{\Phi}^T (\mathbf{y} - \mathbf{t})$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n (1-y_n) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T = \boldsymbol{\Phi}^T \mathbf{R} \boldsymbol{\Phi}, \quad \mathbf{R} = \text{diag}(y_1(1-y_1), \dots, y_n(1-y_n))$$

where in the last equation we used $d\sigma/da = \sigma(1-\sigma)$

- The Hessian now depends on \mathbf{w} through *the weighting matrix \mathbf{R}* , since the error function is no longer quadratic.
- Using $0 < y_n < 1$, which follows from the form of the logistic sigmoid function, we see that $\mathbf{u}^T \mathbf{H} \mathbf{u} > 0$ for any \mathbf{u} , and so **\mathbf{H} is positive definite**.
- It follows that **the error function is a convex function of \mathbf{w} and hence has a unique minimum.**

Iterative Reweighted Least Squares

- The Newton-Raphson update then becomes:

$$\begin{aligned}\mathbf{w}^{(new)} &= \mathbf{w}^{(old)} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) = (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi^T \mathbf{R} \Phi \mathbf{w}^{(old)} - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} z, \text{ where } z = \Phi \mathbf{w}^{(old)} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t})\end{aligned}$$

- The update formula takes the form of a set of normal equations for a weighted least-squares problem.

- Because \mathbf{R} depends on \mathbf{w} , we apply the normal equations iteratively.

- For this reason, the algorithm is known as iterative reweighted least squares - IRLS.

- Rubin, D. B. (1983). [Iteratively reweighted least squares](#). In *Encyclopedia of Statistical Sciences*, Volume 4, pp. 272–275. Wiley.

$\mathbf{w} = \mathbf{0}_D;$
 $w_0 = \log \bar{t} / (1 - \bar{t})$;
repeat

$\eta_i = w_0 + \mathbf{w}^T \mathbf{x}_i$;
 $y_i = \text{sigm}(\eta_i)$;
 $r_i = y_i(1 - y_i)$;
 $z_i = h_i + (t_i - y_i) / r_i$;
 $\mathbf{R} = \text{diag}(r_{1:N})$;
 $\mathbf{w} = \Phi^T \mathbf{R} \Phi^{-1} \Phi^T \mathbf{R} z$;
until converged;

Iterative Reweighted Least Squares

$$\mathbf{w}^{(new)} = (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}, \text{ where } \mathbf{z} = \Phi \mathbf{w}^{(old)} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t}),$$

$$\mathbf{R} = \text{diag} \left(y_1 (1 - y_1), \dots, y_n (1 - y_n) \right)$$

- The elements of the diagonal weighting matrix \mathbf{R} can be interpreted as variances.
- Indeed note that the mean and variance of t in the logistic regression model are given as:

$$\mathbb{E}[t] = 1 \times \sigma(\mathbf{x}) + 0 \times (1 - \sigma(\mathbf{x})) = y$$

$$\text{var}[t] = \mathbb{E}[t^2] - \mathbb{E}[t]^2 = \sigma(\mathbf{x}) - \sigma(\mathbf{x})^2 = y(1 - y)$$

where we used $t^2 = t$ for $t \in \{0, 1\}$.

Iterative Reweighted Least Squares

$$\mathbf{w}^{(new)} = (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}, \text{ where } \mathbf{z} = \Phi \mathbf{w}^{(old)} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t}),$$

$$\mathbf{R} = \text{diag} \left(y_1 (1 - y_1), \dots, y_n (1 - y_n) \right)$$

- We can also interpret *IRLS* as the solution to a linearized problem in the space of the variable $a = \mathbf{w}^T \boldsymbol{\phi}$.
- The quantity z_n , which corresponds to the n^{th} element of \mathbf{z} , can be seen as ***an effective target value*** in this space obtained by making a local linear approximation to the logistic sigmoid function around $\mathbf{w}^{(old)}$.

$$y_n = \sigma(a_n) = \frac{1}{1 + e^{-a_n}}, a_n = \mathbf{w}^T \boldsymbol{\phi}_n, \frac{dy_n}{da_n} = y_n (1 - y_n) \Rightarrow$$

$$t_n \cong y_n + \frac{dy_n}{da_n} \Big|_{\mathbf{w}^{(old)}} (a_n(w) - a_n(w^{(old)}))$$

$$a_n(w) \simeq a_n(w^{(old)}) + \frac{da_n}{dy_n} \Big|_{\mathbf{w}^{(old)}} (t_n - y_n) = \boldsymbol{\phi}_n^T \mathbf{w}^{(old)} - \frac{(y_n - t_n)}{y_n (1 - y_n)} = z_n$$

Multiclass Logistic Regression

- For $K > 2$ classes, the posterior probabilities are given by a *softmax transformation of linear functions of the feature variables* as follows:

$$p(C_k | \phi) = y_k(\phi) = \frac{e^{a_k}}{\sum_j e^{a_j}}, \quad a_k = \ln(p(\phi | C_k) p(C_k)) = \mathbf{w}_k^T \phi$$

- Use maximum likelihood to determine $\{\mathbf{w}_k\}$ for this model.
- We need to compute the derivatives of y_k with respect to all of the activations a_j . They are given as:

$$y_k(\phi) = \frac{e^{a_k}}{\sum_j e^{a_j}} \Rightarrow \frac{\partial y_k}{\partial a_j} = y_k (\delta_{kj} - y_j)$$

- Krishnapuram, B., L. Carin, M. Figueiredo, and A. Hartemink (2005). Learning sparse bayesian classifiers: multi-class formulation, fast algorithms, and generalization bounds. IEEE Transaction on Pattern Analysis and Machine Intelligence.

Multiclass Logistic Regression

- Next we write down the likelihood function.
- We use the 1-of-K coding scheme. The target t_n for a feature ϕ_n belonging to class C_k is a binary vector with all elements zero except for element k , which equals one.
- The likelihood function is then given by

$$p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_{n=1}^N \prod_{k=1}^K p(C_k | \phi_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_{nk}^{t_{nk}},$$

$$y_{nk} = y_k(\phi_n), \mathbf{T} \text{ } N \times K \text{ matrix}, T_{nk} = t_{nk}$$

- The negative log likelihood is then

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

- This is the cross-entropy error function for the multiclass classification.

Multiclass Logistic Regression

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T} \mid \mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

- We now take the gradient of the error function with respect to one of the \mathbf{w}_j .

- Using $y_{nk} = \frac{e^{a_{kn}}}{\sum_j e^{a_{jn}}}$, $\log y_{nk} = a_{kn} - \log \sum_j e^{a_{jn}}$, $a_{kn} = \mathbf{w}_k^T \boldsymbol{\phi}_n$ and $\sum_k t_{nk} = 1$, we obtain

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \boldsymbol{\phi}_n$$

- This is again the form for the gradient derived earlier for the sum-of-squares error function with the linear regression model and the cross-entropy error for logistic regression.
- We *use a sequential algorithm*

$$\mathbf{w}_j^{(\tau+1)} = \mathbf{w}_j^{(\tau)} - \eta \nabla E_{nj} = \mathbf{w}^{(\tau)} - \eta (y_{nj} - t_{nj}) \boldsymbol{\phi}_n$$

General Form: Derivative of log Likelihood

- We have seen that the derivative of the log likelihood function for a linear regression model with respect to the parameter vector \mathbf{w} for a data point n took the form of the ‘error’ $y_n - t_n$ times the feature vector ϕ_n .
- Similarly, for the combination of logistic sigmoid activation function and cross-entropy error function

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1-t_n) \ln (1-y_n) \right\},$$

and for the softmax activation function with the multiclass cross-entropy error function

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = -\ln p(\mathbf{T} | \mathbf{w}_1, \dots, \mathbf{w}_K) = -\sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

- All these are examples of a more general result.

IRLS Algorithm for the Multiclass Problem

- To find a batch algorithm, we again appeal to the Newton-Raphson update to obtain the corresponding IRLS algorithm for the multiclass problem.
- This requires evaluation of the Hessian matrix that comprises blocks of size $M \times M$ in which block j, k is given by

$$\nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N (y_{nj} - t_{nj}) \boldsymbol{\phi}_n \Rightarrow \nabla_{\mathbf{w}_k} \nabla_{\mathbf{w}_j} E(\mathbf{w}_1, \dots, \mathbf{w}_K) = \sum_{n=1}^N y_{nk} (\delta_{kj} - y_{nj}) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T$$

where we used:

$$\nabla_{\mathbf{w}_k} y_{nj} = \frac{dy_{nj}}{da_{nk}} \nabla_{\mathbf{w}_k} a_{nk} = y_{nk} (\delta_{kj} - y_{nj}) \boldsymbol{\phi}_n$$

- As with the two-class problem, the Hessian matrix for the multiclass logistic regression model is positive definite and so the error function again has a unique minimum.

- Bishop, C. M. and I. T. Nabney (2008). *Pattern Recognition and Machine Learning: A Matlab Companion*. Springer. In preparation

MAP Estimate

- One can define an appropriate prior

$$p(\mathbf{W}) = \prod_c \mathcal{N}(\mathbf{w}_c | \mathbf{0}, \mathbf{V}_0)$$

- Using this, the modified negative log-likelihood and its gradient and Hessian are:

$$f'(\mathbf{W}) = -\log p(\mathcal{D}|\mathbf{W}) - \log p(\mathbf{W}) = f(\mathbf{W}) + \frac{1}{2} \sum_c \mathbf{w}_c \mathbf{V}_0^{-1} \mathbf{w}_c$$

$$\mathbf{g}'(\mathbf{W}) = \mathbf{g}(\mathbf{W}) + \mathbf{V}_0^{-1} \sum_c \mathbf{w}_c$$

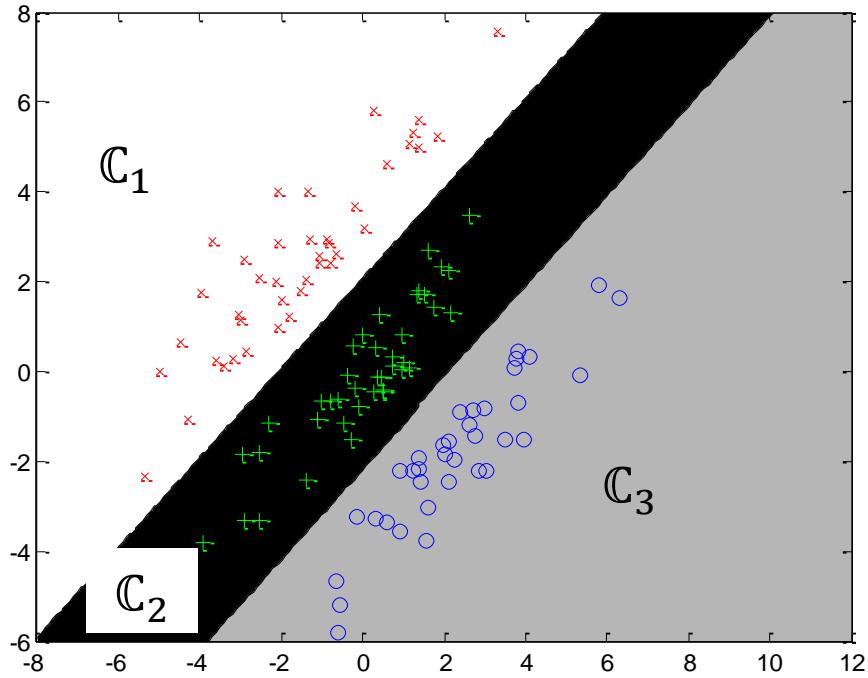
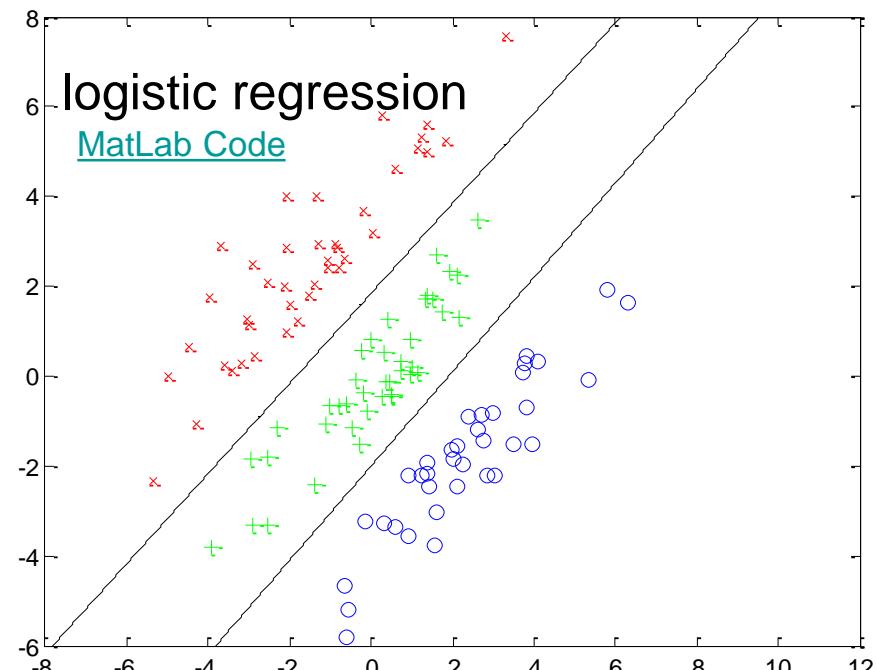
logregFit
from [Kevin Murphys' PMTK](#)

$$\mathbf{H}'(\mathbf{W}) = \mathbf{H}(\mathbf{W}) + \mathbf{I}_C \otimes \mathbf{V}_0^{-1} ((CD) \times (CD))$$

- The symbol $\mathbf{I}_C \otimes \mathbf{V}_0^{-1}$ is used here to denote [Kronecker matrix product](#).
- One can apply [*limited memory BFGS*](#) to find the MAP estimate.

Multiclass Logistic Regression Example

This is the same example of 3 class classification discussed earlier.



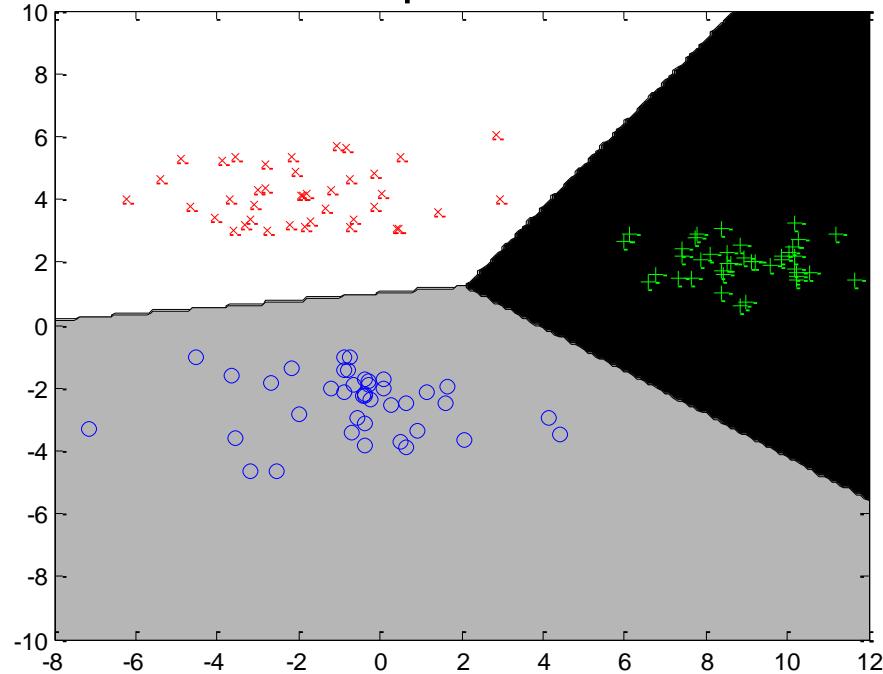
Here, a_1, a_2, a_3 are the feature functions computed by following the logistic regression algorithm. Region C_1, C_2, C_3 in the right plot can be interpreted as:

$$C_m: p(C_m|\phi) = \max(p(C_1|\phi), p(C_2|\phi), p(C_3|\phi)) \text{ where}$$

$$p(C_k|\phi) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

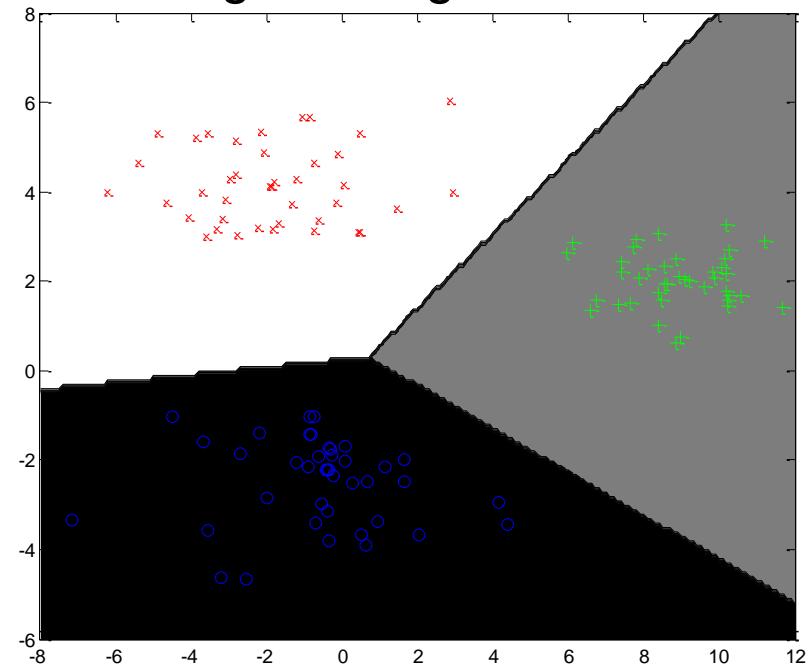
Multiclass Logistic Regression Example

Least-squares



MatLab Code:
Use 'data_3c_2'

Logistic Regression



MatLab Code:
Use 'data_3c_2'